# Applied Statistics:

Books:

1) Johnson, Beaker book.

2) ISLR, Springer.

## Data frames:

$$\begin{cases} n : \text{statistical units.} \\ p : \text{features.} \end{cases}$$

$\underline{x}_1 = (x_{11}, x_{12}, \ldots, x_{1p})' \in \mathbb{R}^p$.

$\underline{x}_2 = (x_{21}, \ldots, x_{2p})'$
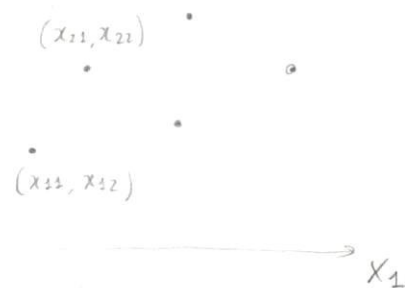
$\vdots$

$\underline{x}_n = (x_{n1}, \ldots, x_{np})'$

Basically, this is a matrix of observations:

$$\begin{array}{c c} & \begin{matrix} X_1 & X_2 & & & X_p \to \text{features} \end{matrix} \\ \text{stat. units:} \begin{matrix} 1 \\ 2 \\ \vdots \\ \vdots \\ n \end{matrix} & \begin{bmatrix} x_{11} & x_{12} & \cdots & & x_{1p} \\ x_{21} & & & & x_{2p} \\ \vdots & & & & \vdots \\ & & & & \\ x_{n1} & \cdots & & & x_{np} \end{bmatrix} = \mathbb{X} \end{array}$$

data matrix (data frame)

$$Y \quad \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

In our usual setting, $\boxed{n \gg p}$

$\underline{X}$ : random vector, $\underline{X} \in \mathbb{R}^p$.

$Y$ : random variable, $Y \in \mathbb{R}$.

$X_2$

$(x_{11}, x_{22})$

$(x_{11}, x_{12})$

$X_1$

## Problem: use $\underline{X}$ to predict $Y$.

We are looking for the "best" function, $f : \mathbb{R}^p \to \mathbb{R}$ to predict $Y$ in terms of $\underline{X}$.

Best in the sense that $\boxed{\mathbb{E}\left[(Y - f(\underline{X}))^2\right]}$ should be minimum.
We are taking the squares because we want to minimize the
mean squared error. We could use different measures, like the absolute
value. For now, we are working in $L^2$.

$$E\left[(Y-f(\underline{x}))^2\right].$$

Exercise:    find that: $\text{Arg}\min\limits_{k} E\left[(Y-k)^2\right] = E[Y]$.
(prove)

In order to solve the more general problem of finding $f$, we do:

$$E\left[(Y-f(\underline{x}))^2\right] = E\left[(Y-E[Y|\underline{x}] + E[Y|\underline{x}] - f(\underline{x}))^2\right] =$$

$$= E\left[(Y-E[Y|\underline{x}])^2\right] + E\left[(E[Y|\underline{x}]-f(\underline{x}))^2\right] +$$

$$+ 2E\left[(Y-E[Y|\underline{x}])(E[Y|\underline{x}]-f(\underline{x}))\right].$$

We have that: given $W, Z,$

$$E[W] = E[E[W|Z]].$$    Hence the middle term becomes:

$$(*) = E\left[(E[Y|\underline{x}]-f(\underline{x}))\, E[Y-E[Y|\underline{x}]\,|\,\underline{x}]\right] \quad \Rightarrow \quad (*) = 0.$$

Condition on $\underline{X}$

$$= E[Y|\underline{x}] - E[Y|\underline{x}] = 0.$$

Then:

$$E\left[(Y-f(\underline{x}))^2\right] = E\left[(Y-E[Y|\underline{x}])^2\right] + E\left[(E[Y|\underline{x}]-f(\underline{x}))^2\right].$$

Thus, our best guess is to take $f(\underline{x}) = E[Y|\underline{x}]$.    solution of the optimization problem

The first term cannot be eliminated; it's like a constant.
There will always be a difference between $Y$ and its prediction $f(\underline{x})$:

$$Y - f(\underline{X}) = \varepsilon$$

This relation is giving us a model:

$$Y = f(\underline{X}) + \varepsilon, \quad \text{where} \quad f(\underline{x}) = E[Y|\underline{x}].$$

Which are the features of $\varepsilon$?

Observe that:

$$E[Y] = E[E[Y|\underline{x}]] + E[\varepsilon]$$
$$= E[Y] + E[\varepsilon] \quad \Rightarrow \quad E[\varepsilon] = 0.$$

We want to use data to estimate $f$.

Say that $\hat{f}$ is my estimate of $f$. I want to see how good is this estimate:

$$\begin{cases} \hat{f} : \text{estimate of } f \ (\text{via } \mathbb{X}). \\ \underline{x}_0 \in \mathbb{R}^p \xrightarrow[\hat{f}]{\text{prediction}} Y_0. \end{cases} \qquad Y_0 = f(\underline{x}_0) + \varepsilon_0.$$

I want to find $\quad E|_{\mathbb{X}}\left[ (Y_0 - \hat{f}(\underline{x}_0))^2 \right] =$
$\underbrace{\qquad}_{\text{conditioning on data}}$

$$= E|_{\mathbb{X}}\left[ (f(\underline{x}_0) + \varepsilon_0 - \hat{f}(\underline{x}_0))^2 \right] =$$

$$= E|_{\mathbb{X}}\left[ \underbrace{(f(\underline{x}_0) - \hat{f}(\underline{x}_0))^2}_{\text{constant}} \right] + \underbrace{E|_{\mathbb{X}}\left[ \varepsilon_0^2 \right]}_{\varepsilon_0 \text{ independent on } \mathbb{X}} + \underbrace{2 E|_{\mathbb{X}}\left[ (f(\underline{x}_0) - \hat{f}(\underline{x}_0)) \varepsilon_0 \right]}_{\substack{\text{constant, can take out} \\ \text{of the expectation}}}$$

$$= (f(\underline{x}_0) - \hat{f}(\underline{x}_0))^2 + Var(\varepsilon_0) + 0.$$

$$\Rightarrow \boxed{E|_{\mathbb{X}}\left[ (Y_0 - \hat{f}(\underline{x}_0))^2 \right] = (f(\underline{x}_0) - \hat{f}(\underline{x}_0))^2 + Var(\varepsilon_0).}$$

$\underbrace{\qquad\qquad}_{\substack{\text{reducible term:} \\ \text{if } \hat{f} \text{ is a good estimate,} \\ \text{it will be small.}}} \qquad \underbrace{\qquad}_{\substack{\text{irreducible} \\ \text{error term}}}$

### Estimation of $f$:



$\hat{f}(\underline{x})$ estimate of $E[Y|\underline{x}]$.
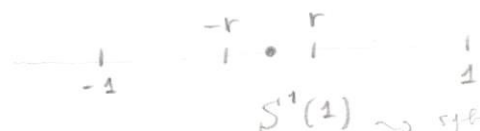
One option is the following:

$\hat{f}(x):$ average of $\{y_i\}$ where $x_i \in N(\bar{x})$.

When $p$ is large, this technique is not any longer feasible $\Rightarrow$ "curse of dimensionality".

### Curse of dimensionality:

p = 1:



$x \sim U[S^1(1)].$

$S^1(1) \rightsquigarrow$ sphere of radius 1 in dimension 1

How big is the distance that you have to travel in order to capture 10% of your friends $x$?

$$0.1 = \frac{\text{length }(S^1(r))}{\text{length }(S^1(1))} = \frac{2r}{2} = r.$$